

The “handedness” of language: Directional symmetry breaking of sign usage in words

Md Izhar Ashraf^{1,2} and Sitabhra Sinha^{1,3,*}

¹The Institute of Mathematical Sciences, CIT Campus, Taramani, Chennai 600113, India.

²B. S. Abdur Rahman University, Seethakathi Estate, Vandalur, Chennai 600048, India.

³National Institute of Advanced Study, Indian Institute of Science Campus, Bangalore 560012, India.

*sitabhra@imsc.res.in

ABSTRACT

Using large written corpora for many different scripts, we show that the occurrence probability distributions of signs at the left and right ends of words have a distinct heterogeneous nature. Characterizing this asymmetry using quantitative inequality measures, we show that the beginning of a word is less restrictive in sign usage than the end. The asymmetry is also seen in undeciphered inscriptions and we use this to infer the direction of writing which agrees with archaeological evidence. Unlike traditional investigations of phonotactic constraints which focus on language-specific patterns, our study reveals a property valid across languages and writing systems. As both language and writing are unique aspects of our species, this universal signature may reflect an innate feature of the human cognitive phenomenon.

Introduction

Language - and by extension, writing - distinguishes humans from all other species.¹ The ability to communicate complex information across both space and time have enabled society and civilization to emerge.² The recent availability of publicly accessible “Big data”, such as the large digitized corpus on the *Google Books* website, has revolutionized the quantitative analysis of socio-cultural phenomena.³ These include questions about human language, such as how vocabularies evolve over time⁴ and the possible existence of universal patterns in the emotional spectrum of languages.⁵ Language in its written form is represented as symbolic sequences that convey information, and is well-suited for statistical analysis to infer underlying patterns. One of the best known empirical regularities associated with language is the scaling behavior - referred to as Zipf’s law - that quantifies how some words occur far more frequently than others.⁶ It is characterized by power-law tails in word frequency distributions (corresponding to the most frequently used words belonging to a “kernel lexicon”⁷), a property that has been validated for frequently used words across many languages, and possible theoretical explanations of the phenomenon have been proposed.^{8,9} Words are themselves composed of letters, and it has long been known that the different letters also occur with characteristic frequencies - a fact that has been used by cryptographers over the ages to break simple substitution ciphers, illustrated dramatically in fiction by Poe (*The Gold-bug*) and Conan Doyle (*The adventure of the dancing men*). For English, the phrase “ETAOIN SHRDLU” has often been used as a mnemonic for recalling the approximate order of the most commonly occurring letters in typical texts. However, a cursory glance through an English dictionary (or encyclopedia) to ascertain, for each letter of the alphabet, the number of pages that are required to list all the words (or entries) that begin with that letter, will alert one to a strong deviation from what is naively expected from the frequency distribution of letters. For instance, one of the letters having the largest number of entries in a dictionary is ‘c’ which does not even appear among the most frequently used letters in English as per the phrase above. This apparent anomaly arises from the fact that the letter ‘c’ has a much higher probability of occurrence (relative to other letters) at the beginning of an English word - possibly a result of the specific orthography of English, where it can appear as the initial letter of the words *china*, *can*, *cent*, etc., in all of which it is pronounced differently - but does not occur so frequently at other positions. While it is rarer to come across situations where words are arranged according to their last character, it is possible to ask whether the frequency distribution of the letters that occur at the end of a word will also similarly show a distinct character.

Fig. 1 (a-c) shows the occurrence probability distribution of the 26 letters of the alphabet at the initial and final positions of words reconstructed from a large database, and compares it with their probabilities to occur anywhere in a text. We note immediately that letters may differ greatly in terms of their occurrence probability depending on the position - but most importantly, the distribution for the right terminal character (the last letter) in an English word appears to be much more heterogeneous than the one corresponding to the left terminal character (the first letter). In other words, the choice of letters that can occur as the final character of a word is more restrictive, i.e., the occurrence probability of the letters are more

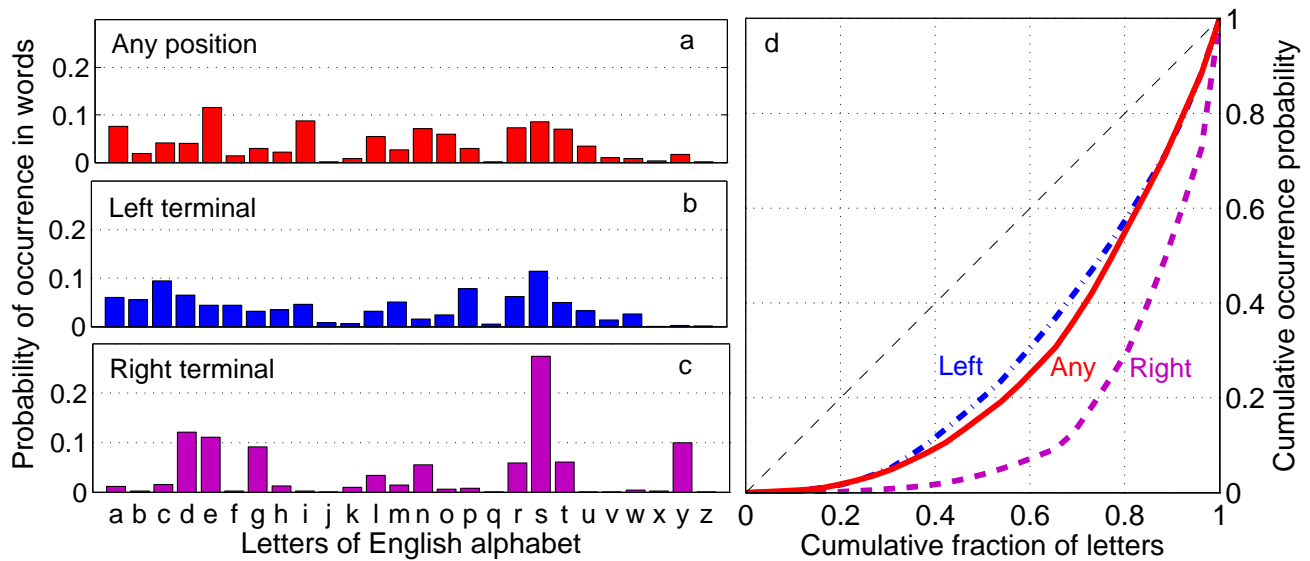


Figure 1. Unequal representation of letters (1-grams) occurring at different positions in words in written English.

The probability of occurrence of the 26 letters of the English alphabet in the *Mieliestronk* corpus comprising about 58000 unique words of the English language (see Methods for details), at (a) any position, (b) left terminal position (i.e., in the beginning) and (c) right terminal position (i.e., at the end) of a word. The distribution shows more heterogeneity for (c), indicating that only a few letters occur with high frequency at the right terminal position of a word, compared to a relatively more egalitarian frequency of occurrence of letters in the left terminal position (b). This difference is illustrated in the Lorenz curve (d) comparing the cumulative distribution function for the occurrence probability of the different letters in any (solid curve), left terminal (dash-dotted curve) and right terminal position (dashed curve) of a word. The thin broken diagonal line corresponds to a perfectly uniform distribution, deviation from which indicates the extent of heterogeneity of letter occurrence probability distributions - measured as the ratio of the area between the line of perfect equality and the observed Lorenz curve, i.e., the Gini index.

unequal, with a very few accounting for the right terminal position for a major fraction of the words, compared to their position-independent probabilities. In contrast, the probability distribution of letters that occur as the initial character is more egalitarian, implying a somewhat higher degree of freedom of choice at the left terminal position. To ensure that this left-right asymmetry in sign usage distributions - suggesting a “handedness” of words in terms of the letter frequency distributions at their terminal positions - is not an artifact of the corpus one is using, we have performed the same analysis with *Google Books Ngram* data, focusing on words that occur with a frequency of more than 10^5 in the corpus digitized by Google. As seen from Fig. S1, the qualitative features are similar to that observed in Fig. 1, although there are differences in the actual frequency of the different letters. This is because, unlike in the preceding case where the occurrence probabilities are computed from a database comprising unique words, the Google 1-gram distributions are computed from corpora of books containing multiple occurrences of the same word - so that the distribution is a joint outcome of the distribution of word frequencies (described by Zipf’s law) and the letter occurrence distribution inferred from a set of unique words. However, despite these differences in details, the inequality of sign usage at the right terminal positions is visibly higher than that in the left terminal positions - indicating the robustness of the observed left-right asymmetry of sign occurrence probability patterns in words.

In this paper we have shown that this directional asymmetry is not just a feature of a particular language but appears to be universal, holding across many languages and writing systems. Regardless of whether the signs we are considering represent letters (for alphabetic scripts like English), syllabograms (for syllabic scripts such as Japanese Kana) or logograms (for logographic scripts like Chinese or logo-syllabic ones like Sumerian cuneiform), the distribution of the signs that begin a word shows relatively less heterogeneity than that for the ones that occur at its end. We have used measures of inequality (such as Gini index or information entropy) to quantitatively assess the degree of asymmetry in the sign occurrence distributions for different linguistic corpora. The difference in the two distributions also indicate the differential information contents of the initial and final characters - and links our result to the statistical and information-theoretic analysis of language.¹⁰ This approach was pioneered by Shannon who used the concept of predictability, i.e., the constraints imposed on a letter by those that have preceded it, to estimate the bounds for the entropy (the amount of information per letter) and redundancy in

English.^{11,12} Considering the consequences of the most prominent structural patterns of texts - viz., the clustering of letters into words - Schürmann and Grassberger subsequently showed that the the average entropy of letters located inside a word are much smaller than that of the letters at the beginning.^{13,14} However, this is true even if one reverses the word - so that terminal letters of words (whether initial or final) have less predictability than those in other positions. Here we ask the relatively simpler question of whether the statistical properties of the left and right terminal characters are different and find a surprising non-trivial asymmetry in the heterogeneity of the respective distributions. Analysis of correlation between sign occurrences in written texts have traditionally focused on the phonotactic constraints of specific languages, e.g., determining the consonants or consonant clusters that are allowed to occur before and after a vowel in any syllable of a given language. While there is considerable variation between different languages as regards the possible arrangements in which consonants and vowels can be combined to make meaningful words, here we show the existence of general patterns that hold across many different language families.

Results

In order to quantify the heterogeneity in sign usage distribution at specific positions in a linguistic sequence, in particular, the beginning or end of a word, we have used the Gini index or coefficient.¹⁵ It measures dispersion in the distribution of a quantity and is widely used in the socio-economic literature to quantify the degree of inequality, e.g., in the distribution of income of individuals or households.¹⁶ The value of the index G expresses the nature of the empirical distribution relative to an uniform distribution, with $G = 0$ if all values of the variable have the same probability of occurrence ("perfect equality") while $G = 1$ corresponds to the extreme situation with the variable always taking up a single value (corresponding to a delta function probability distribution). Thus, if the probability of occurrence of any sign (e.g., the letters 'a-z' in the case of the English alphabet) at the beginning (or end) of a word is about the same, the corresponding Gini index will be close to zero. Otherwise, it is a finite number (≤ 1) whose exact value depends on the extent of inequality in usage of the different signs. Measures related to the Gini index have previously had limited use in the context of linguistic sequences, e.g., to select attributes for decision tree induction in classification for data mining.¹⁷

Using the Gini index on the distributions of letters (1-gram) and pairs of letters (2-gram) that occur at the left and right terminal positions of words in English, we can quantitatively express the visible difference between patterns of unequal occurrence of signs at the two ends seen in Figs. 1 (b-c) and 2 (a-b). Figs. 1 (d) and 2 (d) show the Lorenz curves - a graphical representation of the inequality of a distribution - for the occurrence probability of the different 1- and 2-grams anywhere in a sequence as well as the two terminal positions. Both diagrams clearly show that the probabilities of different signs to occur at the right terminal position, i.e., the end of a word, is more unequal than their occurrence probability in the beginning (i.e., left terminal position of a word), or indeed, anywhere in a sequence. This quantitatively establishes that there is relatively more variation in the letters (or letter pairs) at the start of a word - and conversely less so when ending it. It indicates an inherent left-right asymmetry in the sign usage distribution of words in English that is related to the different degrees of freedom associated with choosing letters that begin and end a word. This asymmetry is more pronounced for letters or 1-grams, as measured by the normalized difference of Gini indices for the two terminal positions, $\Delta G = -0.50$ (for details see Methods), than for 2-grams ($\Delta G = -0.25$) and becomes even less noticeable for higher-order n -grams. We have therefore focused on analysis using 1-grams for the subsequent results reported here.

While this asymmetry in the usage distribution for signs that begin and end words written in English is certainly striking, it would be even more significant if the phenomenon turns out to be valid for linguistic sequences in general. We have, therefore, carried out a systematic investigation of the inequality in sign usage distributions at the terminal positions of sequences that are chosen from languages spanning a broad array of language families (Fig. S2). The writing systems considered are also quite diverse, ranging from alphabetic to logographic, whose corresponding signaries (i.e., the set of distinct characters used for writing in that system) can vary in size from about two dozen to several thousands. Fig. 3 shows the results obtained for the different written corpora we have analyzed, where the degree of asymmetry in sign occurrence at the left and right ends of a sequence is measured by the normalized difference ΔG between the respective Gini indices. The most important feature of our results is the clear distinction that can be made between languages that are conventionally written left to right, such as English, and those which are written right to left, such as Arabic, according to the sign of ΔG obtained for the corresponding corpus. A negative value of ΔG implies that the signs occurring in the left terminal position have relatively more equitable distribution while the sign usage distribution at the right terminal position is more unequal, and conversely for positive ΔG comparatively few signs occur with high frequency at the left end of a sequence than the right end. Thus, our result implies that all languages and writing systems considered here exhibit an asymmetry between the beginning and end of a word in terms of the degree of inequality manifested in their respective sign occurrence probability distributions, the probability in choosing different signs being significantly more heterogeneous at the end than in the beginning.

To ensure that the observed distinction between the sign usage patterns at the two terminal positions of a sequence are significant, we compared our results with those obtained from corpora of randomized sequences, which by design have the

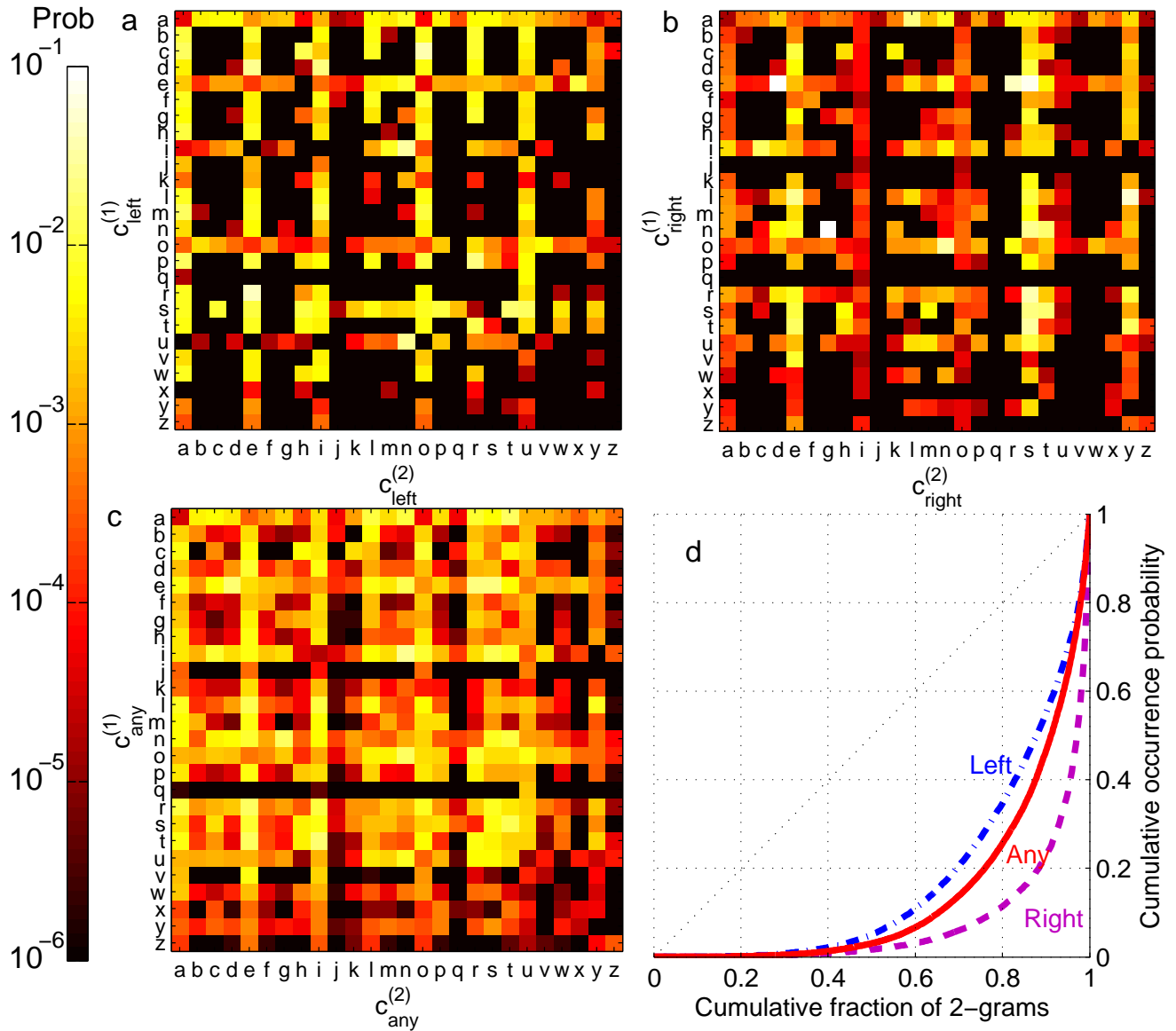


Figure 2. Unequal representation of letter pairs (2-grams) occurring at different positions in words in written English. The probability of occurrence of the 676 ($= 26 \times 26$) different possible letter pairs in the *Mieliestronk* corpus of words (used in Fig. 1) of the English language, at (a) left terminal position (i.e., in the beginning), (b) right terminal position (i.e., at the end) and (c) any position in a word. The corresponding Lorenz curves (d) indicate that, as for 1-grams shown in Fig. 1, the distribution for letter pairs occurring at the right terminal position of a word (dashed curve) show more heterogeneity than those occurring in the left terminal position (dash-dotted curve).

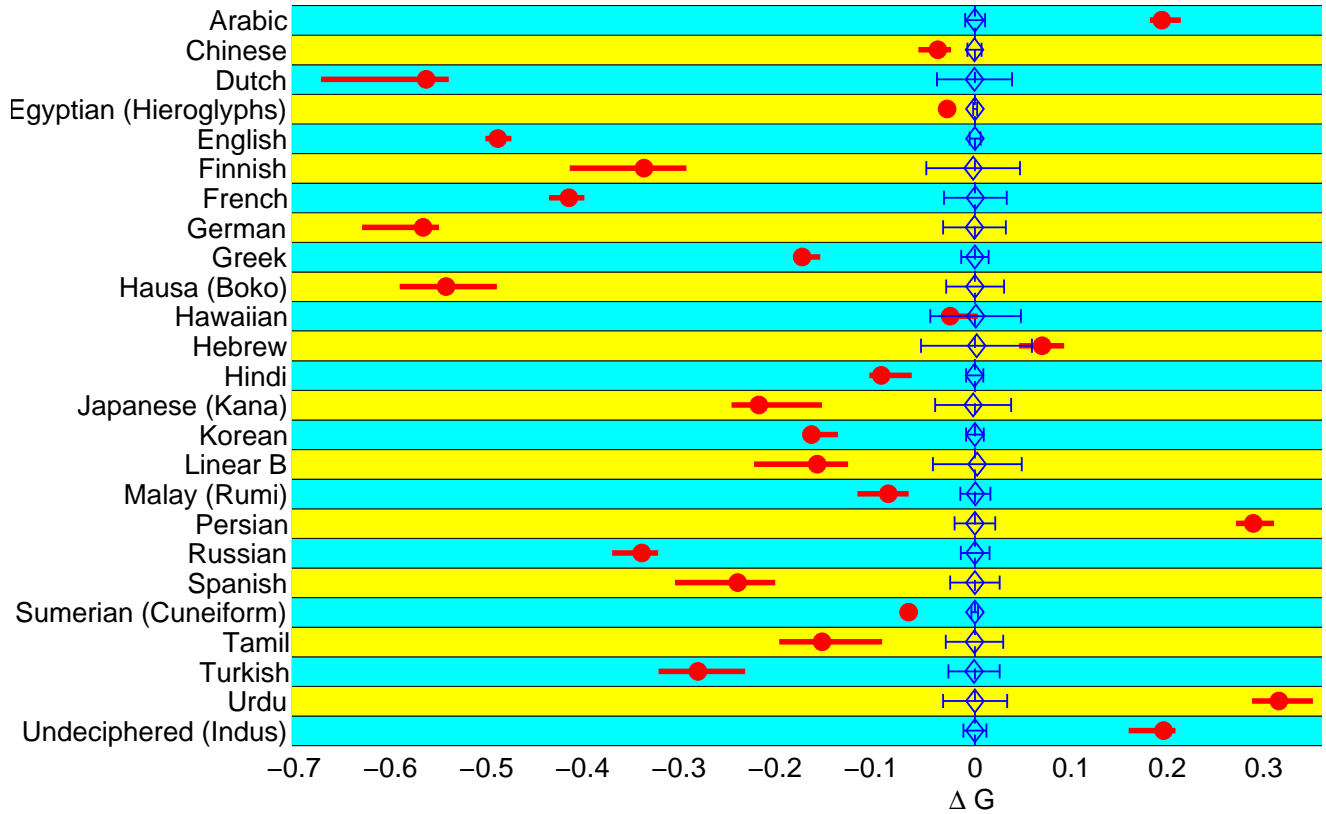


Figure 3. Asymmetry in the sign occurrence probability distributions at the left and right terminal positions of words in different languages correlate with the directions in which they are read. The normalized difference of the Gini indices $\Delta G = 2(G_L - G_R)/(G_L + G_R)$ (filled circles), which measures the relative heterogeneity between the occurrences of different signs in the terminal positions of words of a language, are shown for a number of different written languages (arranged in alphabetical order) that span a variety of possible writing systems - from alphabetic (e.g., English) and syllabic (e.g., Japanese kana) to logographic (Chinese) [see text for details]. All languages that are conventionally read from left to right (or rendered in that format in the databases used here) show a negative value for ΔG , while those read right to left exhibit positive values. The horizontal thick bars superposed on the circles represent the bootstrap confidence interval for the estimated values of ΔG . To verify the significance of the empirical values, they are compared with corresponding ΔG (diamonds) calculated using an ensemble of randomized versions for each of the databases (obtained through multiple realizations of random permutations of the signs occurring in each word). Data points are averages over 1000 random realizations, the ranges of fluctuations being indicated by error bars. Along with the set of known languages, ΔG measured for a corpus of undeciphered inscriptions from the Indus Valley Civilization (2600-1900 BCE) is also shown (bottom row).

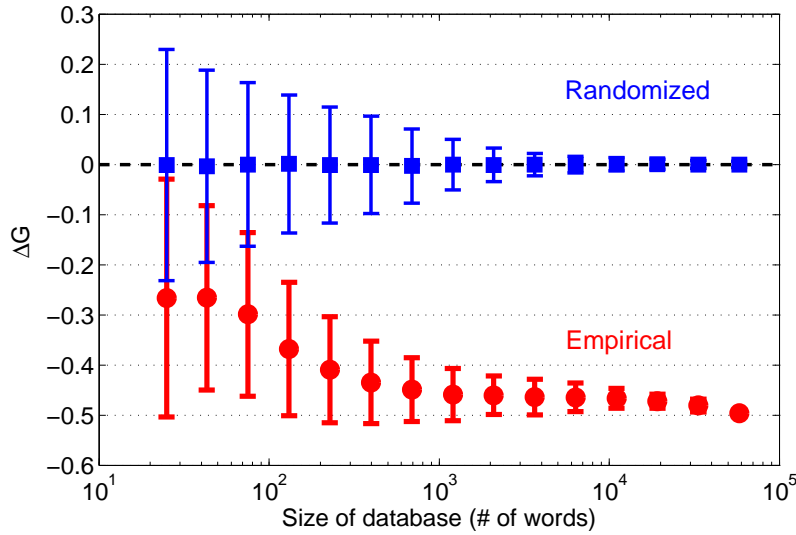


Figure 4. The observed asymmetry between heterogeneity of letter occurrence probability in left and right terminal positions is significant when the database is sufficiently large. Gini index differential ΔG shown for the left and right terminal letter (1-gram) distributions calculated using a set of N words, as a function of N . Empirical results are shown for random samples (without replacement) taken from the *Mieliestronk* corpus comprising about 58000 unique words of the English language, each data point (circles) being the average over 10^3 samples of size N . For each empirical sample, a corresponding randomized sample is created by randomly permuting the letters in each of the N words, and a data point for the randomized set (squares) represents an average over randomizations of 10^3 samples of size N . With increasing N the empirical distribution becomes distinguishable from the randomized set (which, by definition, should not have any left-right asymmetry). The error bars indicate standard deviation over the different samples.

same distribution of sign occurrences at all positions. For a rigorous comparison, we have used surrogate datasets that have the same frequency distribution of different signs as the original corpus (see Methods for details) so that any distinction between them arises only from differences in the nature of the distributions of sign occurrence at the terminal positions. As randomized sequences are expected not to have any left-right asymmetry in sign usage patterns, the mean value of ΔG for the surrogate data is expected to be zero. However, statistical fluctuations will result in the random corpora belonging to the ensemble having small non-zero values of ΔG distributed about 0 and the standard deviation of the distribution (indicated by error bars in Fig. 3) indicates whether a observed difference in Gini indices can arise by chance even when there is no asymmetry. As seen in Fig. 3 almost all the corpora analyzed by us exhibit asymmetry that is clearly distinct from what might be expected if it was only resulting from noise.

As the asymmetry observed in the linguistic sequences should not be sensitively dependent on the particular corpus from which they are chosen, we have obtained confidence intervals for the empirical ΔG values by bootstrap resampling of the data (see Methods for details). Fig. 3 shows that in almost all cases this interval does not have any overlap with the distribution for randomized sequences - indicating that our results are robust with respect to variations in the corpus. The accuracy of the estimate, which is inversely related to the length of the confidence interval, appears to become higher as the database size, i.e., the total number of sequences being considered, is increased. Indeed, Fig. 4 shows that the database needs to be larger than a minimal size (~ 100 words for English) in order for the significance of the observed asymmetry to be established. On using increasingly larger databases, the difference between the empirical and randomized corpora become more pronounced.

Apart from the Gini index, the inequality in sign occurrence distribution at the terminal positions of a sequence can also be measured by other means, e.g., using information or Shannon entropy, a key concept in information theory. Fig. S3 shows that using the normalized difference of entropy ΔS , estimated from the sign occurrence distributions at the left and rights ends of the different linguistic corpora, yields qualitatively similar results to those obtained by using the Gini index. The sign of ΔS in all cases is seen to be consistent with the direction of writing, with left-to-right written languages having positive values of ΔS while those written right-to-left have $\Delta S < 0$. This is in complete agreement with our earlier conclusion that there is relatively more equality in the probability of occurrence of different signs at the beginning of a word than at its end - reflected in the higher non-uniformity for sign usage distribution for the latter. Thus, the asymmetry we observe in linguistic sequences is robust with respect to the specific measure of inequality being used.

Intriguingly, we find that the asymmetry can also appear in a corpus of inscriptions that are so far undeciphered and whose relation to language is therefore not yet established. As an illustration, we have analyzed sign sequences appearing in the archaeological artifacts (e.g., seals, sealings, pottery, copper tablets, etc.) obtained from excavations carried out at sites of the Indus Valley Civilization (IVC) that existed during 2600-1900 BCE in present day Pakistan and northwestern India.¹⁸⁻²⁰ While there is some debate as to whether these inscriptions constitute “writing” in the sense of encoding spoken language,^{21,22} there is unanimity among all serious scholars that these were mostly written from right to left as inferred from the archaeological evidence (e.g., signs get more crowded at the left end of some inscriptions or spill out of an otherwise linear arrangement).²³⁻²⁶ We observe from Fig. 3 that the ΔG for sign usage distribution is positive, indicating that the choice of signs is less restricted in the right terminal position than the left. This would suggest, based on the connection previously seen between the sign of ΔG and the direction of writing, that the IVC inscriptions are written from right-to-left, which corroborates the consensus view as mentioned above.

Discussion

Inferring the direction of writing is one of the basic pre-requisites for interpreting any linguistic sequence. A variety of possible directions have been seen in different writing systems, both historical and present.²⁷ The most common, left to right in horizontal lines, is the direction in which all scripts descending from the Greek and Brāhmī systems are written, including English, French, German, Hindi and Tamil. Scripts that are written in the other direction, i.e., right to left in horizontal lines, are also common and are used in ancient and modern Semitic scripts including Arabic and Hebrew. Another common orientation is from top to bottom in vertical columns, which is the direction in which Chinese and scripts influenced by it (such as Japanese) are traditionally written. Other, less common, directions of writing are also known, including bottom to top (the Celtic Ogham script) and *boustrophedon*, where the direction reverses in successive lines (as in archaic Greek and Luwian hieroglyph inscriptions). In cases where the inscriptions are undeciphered, such as those of IVC, the direction usually has to be inferred by indirect means. The asymmetry in sign usage patterns reported here - which shows that the beginning of sequences can be distinguished from the end by the nature of heterogeneity in the distributions of sign occurrence at these positions - can provide a valuable tool for ascertaining the direction of writing in such cases. Availability of a sufficiently large corpus would however be necessary for a reliable determination of the direction of writing in these inscriptions.

The reason for the appearance of the directional asymmetry in sign usage distributions for linguistic sequences is yet to be definitively identified. However, it is not unreasonable to expect that this is related to the phonotactic constraints inherent in different languages. The initial sound of a word can be chosen with greater freedom from the set of all available speech sounds (phonemes) of the language, compared to all subsequent sounds that may depend - to a greater or lesser extent - on the sound(s) preceding them. For example, very few of the three-consonant clusters that can in principle occur in English are actually allowed.²⁸ Thus, one would expect a higher degree of variability in the initial sound compared to the one at the end of a word. As writing reflects the patterns of spoken language, to greater or lesser extent depending on the system, one would expect this difference between the beginning and end to be manifested in it. An indirect indication that phonotactic considerations may be at least partially responsible for the asymmetry is provided by the degree of the difference between the inequalities of sign usage at the two ends of a word in different writing systems - especially when normalized entropy difference is used as a measure. We observe that, broadly speaking, the magnitude of ΔS (as well as, ΔG) is larger for scripts that have a higher proportion of phonetic representation.^{29,30} Thus, alphabetic and syllabic systems which have a much greater phonetic character than logographic or logo-syllabic systems tend to typically show a more pronounced asymmetry (Fig. S3). As even an apparently logographic system such as Chinese have some degree of phoneticism,²⁹ it is not surprising that systems having a high degree of logography also show a difference in the sign usage distribution between the beginning and end of words, although this effect is much less marked than in other (more phonetic) scripts. There are exceptions from this general trend - for example, Hebrew, which is an alphabetic script, shows a low degree of asymmetry that may be difficult to distinguish from effects due to stochastic fluctuations arising from sampling effects in a finite corpus. Hawaiian, that also appears to have a very low ΔG , however, shows a significant asymmetry when information entropy is used to measure sign usage inequality in place of the Gini index (see supplementary information). Other indications that a simple phonotactic explanation for the observed asymmetry may not be adequate is shown by the fact that the relative position of some languages in terms of ΔG (or ΔS) do not necessarily conform to common perceptions about the degree of phonetic representation in the corresponding scripts used for writing them.³⁰ For example, the Korean han’gŭl script is considered to have a higher proportion of phoneticism than French;³¹ however, the latter exhibits higher asymmetry in terms of both the measures of inequality used here (see Fig.3 and supplementary information).

To summarize, we have reported here evidence for a novel universal feature in the empirical statistics of linguistic sequences. Unlike the more well-known Zipf’s law and Heap’s law, which relate to the frequency of word usage, we focus at a more elementary level, viz., that of the signs - corresponding to letters, syllabograms or logograms, depending on the writing system - which constitute individual words. The distribution of occurrence for the different signs at the left and right terminal

positions in a word are shown to have distinct heterogeneous characters that are characterized by measures of inequality such as the Gini index or information entropy. We observe that, in general, the information content at the beginning of a sequence tends to be higher than at the end, which is reflected in the significant asymmetry in terms of the restriction of sign usage at these the two positions. While traditional linguistic investigation of the patterns of consonants and vowels used for constructing syllables in different languages have tended to focus only on describing which specific combinations are allowed in a particular language, our study of the actual frequency of occurrences of the different signs that appear at specific positions in a sequence result in a richer picture. In particular, it yields a pattern that is valid across different languages and scripts, possibly revealing a feature inherent in the information processing and communicating capabilities of the human cognitive apparatus.

Methods

Data description. We have analyzed data from written corpora of twenty four languages (twenty two belonging to nine linguistic families, as well as, two language isolates), along with a corpus of undeciphered inscriptions from the Indus Valley Civilization (ca. 2600-1900 BCE). The writing systems considered range from alphabetic (that use only a few dozens of distinct letters) and syllabic to logo-syllabic and logographic (involving thousands of signs). The average corpus size is about ten thousand unique words collected from a variety of sources. Each word considered for our analysis consisted of multiple graphemes, corresponding to letters, logograms, hieroglyph signs or syllables depending on the writing system used. Detailed description of each corpus is provided in the supplementary information.

Estimation of occurrence probability distribution. Probability distribution of sign occurrences in a corpus of inscriptions are estimated from frequency counts of the distinct signs appearing in the sequences belonging to the database. For establishing the directional asymmetry of sign usage, we focus specifically on the sign occurrence distributions at the left and right terminal positions of a sequence. The inequality of sign usage at these positions, which is reflected in the non-uniform nature of the corresponding distributions, is quantified by measuring the Gini coefficient or the information entropy.

Measuring Gini coefficient. The Gini coefficient or index is a measure of how unequal are the probabilities of all the different events that are possible, with a value of zero corresponding to the situations where all events are equally probable. Conversely, when only one event out of all possible ones is observed in every instance, the Gini index has the maximum value of 1. For a discrete probability distribution $P(x)$, where the values of the discrete variable x are indexed in increasing order ($x_i < x_{i+1}$, $i = 1, \dots, N$), the Gini coefficient is measured as

$$G = 1 - \sum_{i=1}^N (S_i - S_{i-1})[P_c(x_i) - P_c(x_{i-1})], \quad (1)$$

where $S_i = (\sum_{j=1}^i x_j) / (\sum_{j=1}^N x_j)$ is the cumulative fraction of the variable with $S_0 = 0$ and $S_N = 1$, while $P_c = \sum_{j=1}^i P(x_j)$ is the cumulative probability of x with $P_c(x_0) = 0$ and $P_c(x_N) = 1$. For a given set of inscriptions, we obtain the cumulative probability distributions $P_c^{(L)}$ and $P_c^{(R)}$ for signs occurring in the left and right terminal positions, respectively, and use Eq. 1 to compute the corresponding Gini indices G_L and G_R (considering only the signs which have non-zero probability of occurrence at those positions). The normalized difference between these two values, $\Delta G = 2(G_L - G_R) / (G_L + G_R)$, provides a measure for the asymmetry in the distribution of sign frequencies at the left and right terminal positions of the sequences.

Estimating information entropy. Apart from Gini index, we have used a measure based on information or Shannon entropy for quantifying the nature of the inequality of sign usage distributions at left and right terminal positions in a sequence. As entropy measures the unpredictability of information generated from a source, it can be used to characterize the underlying distribution of any process that produces a discrete sequence of symbols (chosen from a set of N possible ones) and is defined as

$$S = -\sum_{i=1}^N p_i \log_2(p_i), \quad (2)$$

where p_i is the probability of occurrence of the i -th symbol and the use of base 2 logarithm implies that the entropy can be expressed in units of bits.¹¹ In particular, given any database of inscriptions, we obtain the probabilities $p_i^{(L)}$ and $p_i^{(R)}$ for a particular sign i from the corresponding signary to occur in the left terminal and right terminal positions. After obtaining these probabilities for all N signs that occur in the corpus of inscriptions, the left and right terminal entropies (S_L and S_R , respectively) are calculated by using Eq. 2. The normalized difference between the two entropy values, $\Delta S = 2(S_L - S_R) / (S_L + S_R)$, provides a measure of the degree of asymmetry in sign frequency at the two ends of a sequence. For our entropy calculation, we have used all signs that occur in the corpus with a frequency of at least 20% of the mean frequency of occurrence, as including signs that occur extremely rarely can mask the inequality of usage among the relatively common signs.

Bootstrap confidence intervals. To quantify the degree of robustness in the estimates obtained using the empirical databases, we have used a bootstrap method to obtain confidence intervals for the measured values. For each corpus we have created 10^3 resampled datasets (i.e., bootstrap samples) by random sampling with replacement, containing the same number of sequences as the original dataset. The probability distributions for sign usage at the terminal positions are then calculated for every bootstrap sample. Finally, a confidence interval for the normalized difference between the left and right terminal Gini indices, ΔG (or, of information entropy, ΔS) for the corpus is computed using the bias-Corrected and accelerated (BCa) method which adjusts for bias and non-constant variances in the bootstrap sample distributions.³²

Sequence randomization. The statistical significance of the measured asymmetry in sign usage is measured by comparing the results obtained from the empirical database with an ensemble of randomized surrogate sequence corpus. Each ensemble is generated by taking each sequence in turn that belongs to a database and doing a random permutation of the signs. This reordering ensures that the frequency distribution of the signs in each sequence (and thus, also the corpus) is unchanged in the randomized set, although all correlations (that contribute to 2- and higher order n -gram distributions) are disrupted. We then perform the same calculations as for the original empirical data for measuring the degree of asymmetric sign usage in left and right terminal positions of these randomized sequences - which, by design, are expected not to have any asymmetry. Significant difference in the results of the two datasets ensures that the measured asymmetry is not arising from stochastic fluctuations.

References

1. Deacon, T. W. *The symbolic species: The co-evolution of language and the brain* (W. W. Norton, New York, N.Y., 1997)
2. Dunbar, R. *Grooming, gossip, and the evolution of language* (Faber, London, 1996).
3. Michel, J.-B. *et al.* Quantitative analysis of culture using millions of digitized books. *Science* **331**, 176-182 (2011).
4. Petersen, A. M., Tenenbaum, J. N., Havlin, S., Stanley, H. E. & Perc, M. Languages cool as they expand: Allometric scaling and the decreasing need for new words. *Scientific Reports* **2**, 943 (2012).
5. Dodds, P. S. *et al.* Human language reveals a universal positivity bias. *Proc. Natl. Acad. Sci. USA* **112**, 2389-2394 (2015).
6. Zipf, G. *Selected studies of the principle of relative frequency in language* (Harvard University Press, Cambridge, Mass., 1932).
7. Cancho, R. F. i & Solé, R. V. The small world of human language. *Proc. Roy. Soc. Lond. B* **268**, 2261-2265 (2001).
8. Mitzenmacher, M. A brief history of generative models for power law and lognormal distributions. *Internet Mathematics* **1**, 226-251 (2003).
9. Ferrer i Cancho, R. & Solé, R. V. Least effort and the origins of scaling in human language. *Proc. Natl. Acad. Sci. USA* **100**, 788-791 (2003).
10. Mumford, D. & Desolneux, A. *Pattern Theory: The stochastic analysis of real-world signals* (A. K. Peters, Natick, Mass., 2010).
11. Shannon, C. E. A mathematical theory of communication. *Bell System Tech. J.* **27**, 379-423 (1948).
12. Shannon, C. E. Prediction and entropy of printed English. *Bell System Technical Journal* **30**, 50-64 (1951).
13. Schürmann, T. & Grassberger, P. The predictability of letters in written English. *Fractals* **4**, 1-5 (1996).
14. Schürmann, T. & Grassberger, P. Entropy estimation of symbol sequences. *Chaos* **6**, 414-427 (1996).
15. Gini C. *Variabilità e Mutuabilità: Contributo allo studio delle distribuzioni e delle relazioni statistiche* (C. Cuppini, Bologna, 1912) [Eng. trans. of extracts in Ceriani, L. & Verme P. The origins of the Gini index. *J. Econ. Inequal.* **10**, 421-433 (2012)].
16. Sinha, S., Chatterjee, A., Chakraborti, A. & Chakrabarti, B. K. *Econophysics: An Introduction* (Wiley-VCH, Weinheim, 2011).
17. Breiman, L., Friedman, J. H., Olshen, R. A., & Stone, C. J. *Classification and Regression Trees* (Wadsworth, Belmont, CA, 1984).
18. Possehl, G. *The Indus Civilization: A Contemporary Perspective* (AltaMira Press, Lanham, MD, 2002).
19. Sinha, S., Ashraf, M. I., Pan, R. K. & Wells, B. K. Network analysis of a corpus of undeciphered Indus civilization inscriptions indicates syntactic organization. *Computer Speech and Language* **25**, 639-654 (2011).
20. Wells, B. K. *The Archaeology and Epigraphy of Indus Writing* (Archaeopress, Oxford, 2015).

21. Lawler, A. The Indus Script - Write or Wrong ? *Science* **306**, 2026-2029 (2004).
22. Rao, R., Yadav, N., Vahia, M., Joglekar, H., Adhikari, R. & Mahadevan, I. Entropic evidence for linguistic structure in the Indus script. *Science* **324**, 1165 (2009).
23. Hunter, G. *The Script of Harappa and Mohenjodaro and its connection with other scripts* (Kegan Paul, London, 1934).
24. Mahadevan, I. *The Indus Script: Texts, Concordance and Tables* (Archaeological Survey of India, Calcutta, 1977).
25. Parpola, A. *Deciphering the Indus Script* (Cambridge University Press, New York, 1994).
26. Daniels, P. T. & Bright W., Eds., *The World's Writing Systems* (Oxford University Press, New York, 1996).
27. Coulmas, F. *The Blackwell encyclopedia of writing systems* (Blackwell, Malden, MA, 1996).
28. McMahon, A. *An Introduction to English Phonology* (Edinburgh University Press, Edinburgh, 2002).
29. DeFrancis, J. & Marshall Unger, J. Rejoinder to Geoffrey Sampson, "Chinese script and the diversity of writing systems". *Linguistics* **32**, 549-554 (1994).
30. Robinson, A. *Writing and Script* (Oxford University Press, Oxford, 2009).
31. Marshall Unger, J. *Ideogram: Chinese Characters and the Myth of Disembodied Meaning* (University of Hawai'i Press, Honolulu, 2004).
32. Efron, B. (1987) Better bootstrap confidence intervals. *J. Amer. Stat. Asso.* **82**, 171-185 (1987).

Acknowledgements

This work was partially supported by the IMSc PRISM project funded by the Department of Atomic Energy, Government of India. We thank P. P. Divakaran, Deepak Dhar, Iravatham Mahadevan, Shakti N. Menon, Adwait Mevada and Chandrasekhar Subramanian for helpful discussions and suggestions. We also thank Bryan K. Wells for allowing the use of the IVC sequence database compiled by him.

Supplementary Information

Description of the corpora

Arabic: We have used a database of 14867 unique words (that are represented using two or more characters) of Classical (or Quranic) Arabic, a Semitic language written using a consonantal alphabet or ‘abjad’.¹ The words are obtained from *Tanzil*, an international project started in 2007 to produce a standard Unicode text for the Qur’an (<http://tanzil.net/download/>, accessed: 25th March 2015). The signary comprises 36 signs, viz., 28 consonantal signs and 8 consonants with diacritical marks indicating vowels. The words range in length from 2 to 11 characters, the average length being 5.39.

Chinese: We have used a database of 10332 unique words and phrases in both traditional and simplified Chinese (belonging to the Sino-Tibetan language family) represented using multiple signs, obtained from a public-domain online *Chinese-English dictionary* CC-CEDICT (<http://www.mdbg.net/chindict/chindict.php?pahe=cc-cedict>, accessed: 8th January 2011). The Chinese writing system is logographic¹ and the signary for our data comprises 3303 distinct graphemes corresponding to logograms (*hanzi*). The sign sequences range in length from 2 to 16 signs, the average length being 5.57. Traditionally, Chinese is written top to bottom in vertical columns shifting from right to left; however, in modern times it is more frequently being written left to right in horizontal lines, and this is the convention used in our database.

Dutch: We have used a list of the 10000 most commonly used words in Dutch, a member of the Germanic branch of the Indo-European language family, from which we have chosen the 9895 words that have two or more characters. The data has been collected from the *Wortschatz* website maintained by the University of Leipzig (<http://wortschatz.uni-leipzig.de/Papers/top10000nl.txt>, accessed: 22nd May 2015). The signary used has 32 distinct alphabetic characters comprising 21 consonants, 5 vowels, 3 vowels with diacritical marks (acute accents or diaeresis), the digraph ‘ij’ that is considered as a letter in the Dutch language, an extra letter from the German alphabet (the *Eszett*) and an apostrophe sign. The words range in length from 2 to 26 letters, the average length being 7.5.

Egyptian (Hieroglyphs): Ancient Egyptian, a member of the Afro-Asiatic (Hamito-Semitic) language family, is written using a mixed system (also referred to as a logoconsonantal system²) with several hundreds of *hieroglyph* signs that can represent logograms, phonograms and/or determinatives.¹ We have used as data 39934 entries (comprising two or more hieroglyph signs) of the *Middle Egyptian Dictionary* compiled by Mark Vygus (updated April 2015, <http://www.pyramidtextsonline.com/MarkVygusDictionary.pdf>, accessed: 22nd May, 2015). The hieroglyphic signs are represented using the Gardiner sign list numbering system,³ the signary for the database used by us comprising 1860 distinct signs. The sign sequences range in length from 2 to 17 hieroglyphs, the average length being 5.11. The conventional direction of reading hieroglyphic sequences is “toward the face of human or animal pictograms, i.e., the signs are turned towards the beginning of the inscription”.² In the database used by us all sequences have been oriented so as to read from left to right.

English: We have used the *Mieliestronk* list of 58111 distinct words (comprising two or more letters) of the English language - belonging to the Germanic branch of the Indo-European language family - that has been compiled by merging several different word-lists (<http://www.mieliestronk.com/wordlist.html>, accessed: 4th December 2011). The words vary in length from 2 to 22 letters, the average being 8.34. The signary consists of the 26 lower case letters of the English alphabet. The list excludes spellings that are considered to be non-British. If a word is hyphenated, it is listed in unhyphenated form by removing the punctuation mark. The list contains some multiword phrases that are in common usage, rendered as a single word. Several words are included in both their singular and plural forms.

For corroboration of our results obtained by analyzing the above dataset, we have also used statistics of sequence position-specific distributions of letter usage frequencies in a list of 97565 distinct words of the English language compiled from *Google books Ngram* data (English Version 20120701, in <http://storage.googleapis.com/books/ngrams/books/datasetsv2.html>) by Peter Norvig and made freely available for public access (<http://norvig.com/tsv/ngrams-all.tsv.zip>, accessed: 22nd June 2015).⁴ Only those words are used which occur with a frequency of more than 100,000 in the corpus of books scanned by Google. Note that, unlike the other data used by us, the frequencies of the different N-grams are weighted by the number of occurrences of each word in the corpus. Details of the procedure used for compiling the N-gram frequency statistics are given in <http://norvig.com/mayzner.html> (accessed: 22nd June 2015).

Finnish: We have used a list of the 10000 most commonly used words (all of which use two or more letters) in the Finnish language, belonging to the Finnic branch of the Uralic language family. The data, obtained from the *Wikiverb* website, has been collected from newsgroup discussions, press and modern literature (<http://wiki.verbix.com/Documents/WordfrequencyFi>, accessed: 24th June 2015). The signary used has 25 distinct signs - i.e., all vowels and consonants of the modern Latin alphabet excepting “q”, “x” and “w”, along with two additional vowels “ä” and “ö”. The words vary in length from 2 to 25 letters, the average being 7.81.

French: We have used a list of the 10000 most commonly used words in French, a Romance language belonging to the Indo-

European family, from which we have chosen the 9826 words that have two or more characters. The data has been collected from the *Wortschatz* website maintained by the University of Leipzig (<http://wortschatz.uni-leipzig.de/Papers/top10000fr.txt>, accessed: May 22nd 2015). The signary used has 30 distinct alphabetic characters comprising 26 letters of the Latin alphabet along with 3 vowels with diacritical marks (acute accents or diaeresis) and an apostrophe sign. The words range in length from 2 to 19 letters, the average length being 7.6.

German: We have used a list of the 10084 most commonly used words in German, a member of the Germanic branch of the Indo-European language family, from which we have chosen the 10053 words that have two or more characters. The data has been collected from the *Wortschatz* website maintained by the University of Leipzig (<http://wortschatz.uni-leipzig.de/Papers/top10000de.txt>, accessed: May 22nd 2015). The signary used has 32 distinct alphabetic characters comprising the 26 letters of the Latin alphabet along with 4 vowels having diacritical marks (umlauts or acute accents), a ligature (the *Eszett* or *scharfes S*) and an apostrophe sign. The words vary in length between 2 to 27 letters, the average being 7.9.

Greek: We have used a list of the 10000 most frequently occurring words - grouped by lemma - in classical Greek literature written in ancient Greek which belongs to the Indo-European language family, compiled by Kyle Johnson from the *Thesaurus Linguae Graecae* corpus (maintained by University of California, Irvine) using the *Classical Language Toolkit* (<http://cltk.org>) and made freely available to the public (<http://kyle-p-johnson.com/assets/most-common-greek-words.txt>, accessed: 24th June 2015). From this dataset we have used the 9889 words that have two or more characters. The signary has 124 distinct characters as the words are represented in the traditional polytonic orthography used for ancient Greek, involving 24 basic letters used in conjunction with several varieties of diacritical marks (e.g., accents, breathing marks, iota subscript and diaeresis). The words range in length from 2 to 18 characters, the average being 6.9.

Hausa (Boko): Hausa, a Chadic language belonging to the Afro-Asiatic family, is written using Boko, a Latin-based alphabet, which was devised in the 19th century and became the official system in the early part of the 20th century (in earlier periods, it was written in Ajami, an Arabic alphabet). We have used a list of 7066 words that have two or more characters obtained from a Hausa online dictionary maintained by the University of Vienna (<http://www.univie.ac.at/Hausa/KamusTDC/CD-ROMHausa/KamusTDC/ARBEIT2.txt>, accessed: 19th May, 2015). The signary used has 30 distinct alphabetic characters comprising 23 letters from the Latin alphabet, four additional signs representing glottalized consonants, two digraphs ('sh' and 'ts') and an apostrophe sign. The words range in length from 2 to 22 characters, the average being 6.0.

Hawaiian: Hawaiian is Polynesian language belonging to the Austronesian family that had no written form until the 19th century when foreign missionaries devised an alphabetic system for recording it based on the Latin script. The data used for our analysis has been collected from the entries of *A dictionary of the Hawaiian language* (1922) compiled by Lorrin Andrews and revised by Henry H Parker (Board of Commissioners of Public Archives of the Territory of Hawaii, Honolulu) and freely available online (<http://ulukau.org/elib/cgi-bin/library?c=parker&l=en>, accessed on 28th May 2015). After removing all non-native words that contain characters that do not belong to the Hawaiian alphabet, we have compiled a data-base of 17588 unique words containing two or more characters. The signary comprises 28 distinct characters, with 12 basic letters - representing 5 vowels and 7 consonants - of the Hawaiian alphabet along with vowels used in conjunction with diacritical marks (breve and macron) indicating short or long pronunciation, and a sign to indicate glottal stop (the '*okina*'). The words range in length from 2 to 26 characters, the average being 8.72.

Hebrew: We have used a list of the 10000 most commonly used words (compiled from online written texts) in modern Hebrew, a Semitic language written using a consonantal alphabet or 'abjad', from which we have chosen the 9994 words that are represented using two or more characters. The data has been collected from a list maintained by *Teach Me Hebrew*, an online Hebrew language learning site (<http://www.teachmehebrew.com/hebrew-frequency-list.html>, accessed: 26th December 2013). The signary comprises 31 signs, viz., 27 consonantal signs (comprising 22 letters of which five use different forms - called *sofit* - when used at the end of a word) and 4 signs used in conjunction with *niqqud* diacritical marks. The words range in length from 2 to 13 characters, the average length being 5.08.

Hindi: Hindi is an Indo-Aryan language, a branch of the Indo-European family, which is written in the Devanagari script that is sometimes classified as an alphasyllabary² or 'abugida'.⁵ Like the other writing systems that are descended from the Brāhmī script of ancient India, Devanagari uses as its main functional unit the *aksara*, which may consist of only a vowel but more frequently represents a syllable consisting of a consonant and an inherent vowel along with diacritical marks that may indicate use of other vowels.¹ We have used a database of 6443 words written using two or more characters, collected from an online dictionary (*Shabdanjali*) of Hindi developed by the Language Technology Research Center at Indian Institute of Information Technology, Hyderabad (<http://ltrc.iiit.ac.in/showfile.php?filename=downloads/shabdanjali-stardict/index.html>, accessed: 20th May 2015). The signary comprises 561 distinct signs, comprising 11 vowels, 33 consonants, their conjunctions with each other and with consonantal sound modifiers (the *anusvara*, *chandrabindu*, *visarga* and *halant*). The words range in length from 2 to 10 signs, the average length being 3.5.

Japanese (Kana): Japanese, which belongs to the Japonic language family, is written using a combination of the logographic *Kanji* system (adopted from Chinese characters) and the syllabic *kana* system. The latter, in turn, consists of a pair of distinct

syllabaries: *hiragana*, used for writing native Japanese words and *katakana*, which is used for foreign words. For our study, we have focused only on the syllabic writing system for Japanese. We have used a list of 1232 words written using two more signs from the kana syllabary, which is obtained from a list of common Japanese words collected from textbooks used by foreign learners of the language and maintained by *Japanese Words*, an online site for learning the Japanese language (<http://www.japanesewords.net/36/over-1000-japanese-words-list/>, accessed: 29th May 2015). The signary has 103 distinct characters comprising 46 basic signs of Hiragana and 21 basic signs of Katakana, 22 Hiragana and 9 Katakana signs used in conjunction with diacritical marks (the *dakuten* and *handakuten*), smaller forms of 4 hiragana characters (viz., of *ya*, *yu* and *yo* which indicate the *yōon* feature, and the *sokuon* used to mark a geminate consonant) and a special symbol (*chōonpu*, the long vowel mark). The words range in length from 2 to 13 characters, the average being 3.8.

Korean: Korean is a language isolate with no established connection to any of the major language families of the world and is written using Han'gŭl, a purely phonetic script, although in earlier times a system based on Chinese characters (Hanja) was used. Each character corresponds to a syllable, the syllabic block being composed of two to six letters (including at least one consonant and one vowel) from the basic alphabet comprising 10 vowels, 14 consonants and 27 digraphs. The number of possible distinct syllabic blocks or characters exceeds 11,000 although a far smaller number is in actual use.¹ We have used a list of 5888 commonly used words compiled by the National Institute of Korean Language in 2004 (publicly accessible from https://en.wiktionary.org/wiki/Wiktionary:Frequency_lists/Korean_5800, accessed: 24th June 2015) from which we have chosen 5406 words that are represented using multiple syllabic blocks. The signary comprises 923 distinct characters (each corresponding to a syllabic block). The words range in length from 2 to 6 characters, the average length being 2.7.

Linear B: Linear B is a syllabic script, with most of its signs representing consonant-vowel combinations, that was used for writing archaic Greek between 1500 and 1200 BCE. We have used as data 1933 entries (comprising two or more characters) of the *Linear B Lexicon* compiled by Chris Tselentis (<https://www.scribd.com/doc/56265843/Linear-B-Lexicon>, accessed: 15th May 2015). The signary comprises 87 distinct signs representing syllables. The words range in length from 2 to 8 signs, the average being 3.8.

Malay (Rumi): We have used a list of the 10000 most commonly used words in Malay, a member of the Austronesian language family, from which we have chosen the 9970 words that have two or more characters. All words are written in *Rumi* or Latin script, which is the most commonly used form for writing Malay at present, although a modified Arabic script (*Jawi*) also exists. The data has been collected from the list of high frequency words that are publicly available at *Invoke IT Blog* (<https://invokeit.wordpress.com/frequency-word-lists/>, accessed: 4th January, 2014). The signary comprises the 26 letters of the Latin alphabet. The words range in length from 2 to 17 letters, the average being 6.8.

Persian: We have used a list of 10000 most commonly used words (represented using two or more characters) in Persian, a member of the Indo-Iranian branch of the Indo-European language family, which is written using a modified form of the consonant Arabic alphabet or 'abjad'. The words are obtained from a list of high-frequency words compiled using the Tehran University for Persian Language corpus and available at *Invoke IT Blog* (<https://invokeit.wordpress.com/frequency-word-lists/>, accessed: 4th January 2014). The signary comprises 40 signs, viz., 32 consonantal signs, a long vowel indicator ('alef madde'), a ligature ('lām alef'), a diacritic ('hamze'), 3 consonants with the 'hamze' diacritical mark and different forms for the consonants 'kāf' and 'ye' when they occur in final position. The words range in length from 2 to 13 letters, the average being 5.2.

Russian: We have used a list of 9103 words that use two or more characters in Russian, a member of the Slavic branch of the Indo-European language family and which is written using a Cyrillic alphabet. The data has been collected from *Russian Learners' Dictionary: 10,000 words in frequency order* compiled by Nicholas J Brown (Routledge, London, 1996), after removing all words that use characters not in the standard Russian alphabet. The signary comprises the 33 letters of the modern Russian alphabet. The words range in length from 2 to 21 letters, the average being 8.0.

Spanish: We have used a list of 5104 high-frequency words (that use two or more characters) in Spanish, a Romance language belonging to the Indo-European family. The data has been collected from *A Frequency Dictionary of Spanish* compiled by Mark Davies (Routledge, London, 2006). The signary used has 35 distinct alphabetic characters comprising 26 letters of the basic Latin alphabet along with an additional character ñ and two digraphs ('ch' and 'll'), as well as, vowels with diacritical marks (acute accents or dieresis). The words range in length from 2 to 19 letters, the average being 7.4.

Tamil: Tamil is a Dravidian language is written in a script (sometimes classified as an 'abugida') derived from Brāhmī script and thus shares a common origin with the Devanagari script used for writing Hindi (see above) although it differs significantly both in appearance and structure.¹ It has 31 basic signs consisting of 12 vowels, 18 consonants and a special character, with combinations of the different vowels and consonants yielding a possible 216 compound letters. Additional characters from the Grantha script and diacritical marks are sometimes used to represent sounds not native to Tamil, e.g., in words borrowed from other languages. As with other Indian scripts, Tamil uses the *aksara* as its basic unit - however, unlike then it has eliminated most conjuncts, consonant clusters being placed in a linear string. The data used for our analysis has been collected from texts (e.g., Paripaadal, Thiruppavai, Kamba Ramayanam, Sundara Kandam, Akananooru songs, etc.) available in *Chennai*

Library, an online repository of Tamil literature (<http://www.chennaiLibrary.com>, accessed: 27th December 2012). From this a data-base of 1991 unique words containing two or more characters was compiled. The signary comprises 187 distinct signs corresponding to different basic and compound letters and the special character. The words range in length from 2 to 9 letters, the average being 3.8.

Turkish: We have used a list of 10000 high-frequency words (that use two or more characters) in Turkish, a member of the Turkic language family. The data has been collected from a *Wiktionary* word frequency list (https://en.wiktionary.org/wiki/Wiktionary:Frequency_lists/Turkish_WordList_10K, accessed: 14th July 2015). The signary used has 32 letters, comprising 29 letters of the Turkish alphabet and 3 vowels used in conjunction with circumflex accents. The words range in length from 2 to 17 letters, the average being 6.9.

Sumerian (Cuneiform): Sumerian, a language isolate spoken in ancient Mesopotamia during the 3rd millennium BCE, was written using a logosyllabic system with several hundred signs of a cuneiform script representing logograms, phonograms and/or determinatives. We have used as data 19221 unique words (comprising two or more cuneiform signs) collected from texts available in the *Electronic Text Corpus of Sumerian Literature* (ETCSL, <http://etcsl.orinst.ox.ac.uk>, accessed: 20th October 2010). The signary comprises 1364 distinct transliteration values of Cuneiform signs that represent syllables. The sequences range in length from 2 to 11 signs, with the average being 3.0.

Urdu: We have used a database of 4998 unique words (that are represented using two or more characters) in Urdu, an Indo-Aryan language belonging to the Indo-European family, that is written using an extended Persian alphabet. The words are obtained from a list of frequently used words maintained by the Center for Language Engineering at Lahore (http://www.cle.org.pk/software/ling_resources/UrduHighFreqWords.htm, accessed: 1st January 2014). The signary comprises 46 signs, viz., 38 consonantal signs, 3 long vowels ('alef madde', 'lām alef madde' and 'ya'), 2 semi-consonants ('hamzah' used in conjunction with 'wao' or 'ya'), a nasalized consonant ('noon ghunna'), a ligature ('lām alef') and an additional sign ('ta' marbuta') used for writing certain loan-words. The words range in length from 2 to 11 letters, the average being 4.6.

Undeciphered (Indus): As an example of an undeciphered corpus on which to apply our analysis, we have used the set of inscriptions obtained from archaeological excavations at various sites of the Indus Valley civilization (ca. 2600-1900 BCE). The data used for our analysis is collected from the comprehensive data-base compiled by Bryan K. Wells^{6,7} from which we have removed all incomplete and multiple-line inscriptions thereby obtaining 1752 unique sequences that contain two or more signs. The Indus signs are represented using the Wells sign list numbering system,^{6,7} the signary for the database used by us comprising 572 distinct signs. The sequences range in length from 2 to 13 signs, the average being 4.6. The direction of the sign sequences vary, the majority being written right to left, although examples of left to right or *boustrophedon* also exist.⁸ In the database used by us, all sequences have been oriented so as to read from right to left.

Robustness of results

An important consideration when quantifying the inequality of sign usage is the size of the signary, i.e., the number of visually distinct signs (sometimes referred to as 'graphs'⁹) that can be identified in each corpus. In several scripts, complex characters that are recognizably the compound of two or more basic characters are quite commonly used - as in the system of conjunct consonant signs or ligatures in the Devanagari script used for writing Hindi and other South Asian languages - and in principle, one could either consider these as separate signs or decompose them into the constituent signs, which will result in very different signary sizes. Also, Semitic scripts such as Arabic and Hebrew that are essentially consonantal alphabets often use diacritical marks for indicating vowel usage. In such cases, the same consonant is used in conjunction with different diacritics when the vowel following it is different. The signary size would depend upon whether these are considered to be distinct signs or not. In several other scripts, special marks can be used together with the vowels and consonants, e.g., the use of apostrophe to indicate the omission of one or more sounds in European languages such as French or German, and the use of a glottal stop marker ('*okina*') in Hawaiian. Once again, whether these signs are treated as distinct elements or part of the associated letter will affect the signary size. We observe that although the numerical values of the Gini indices (and information entropy) can be affected by changes in the signary size, the asymmetry in terminal sign usage reported here is robust with respect to these choices about different conventions for identifying distinct signs constituting the signary for a given corpus.

Possible non-phonotactic mechanisms for emergence of asymmetry

It is possible that the asymmetry we have reported here could arise in certain situations for reasons other than phonotactic constraints. For instance, in the undeciphered IVC inscriptions, the most frequently occurring sign - viz., the U-shaped "jar"

symbol¹⁰ - that appears very frequently at the end of a sequence, dominates the probability distribution of signs that can occur at the left terminal position (accounting for about a third of all the distinct sequences comprising the corpus). By contrast, the most frequently occurring sign at the right terminal position is seen to begin only about 6 % of the sequences. These distinctive sign usage patterns at the two terminal positions of IVC sequences gives rise to the heterogeneity in the corresponding occurrence probability distributions. In the absence of a decipherment, it is purely speculative whether the inequality arises for phonotactic reasons (as in the linguistic sequences considered here) or a fortuitous stylistic convention.

References

1. Coulmas, F. *Writing Systems: An Introduction to their Linguistic Analysis* (Cambridge University Press, Cambridge, 2003).
2. Daniels, P. T. and Bright, W., Eds. *The World's Writing Systems* (Oxford University Press, New York, 1996).
3. Gardiner, A. *Egyptian Grammar: Being an Introduction to the Study of Hieroglyphs* (3rd edn., Griffith Institute, Oxford, 1957).
4. Norvig, P., *English letter frequency counts: Mayzner Revisited or ETAOIN SRHLDCU*. (2013) Available at: <http://norvig.com/mayzner.html>. (Accessed: 22nd June 2015)
5. Daniels, P. T. Fundamentals of grammatology. *J. Amer. Oriental Soc.* **100**, 727-731 (1990).
6. Sinha, S., Ashraf, M. I., Pan, R. K. & Wells, B. K. Network analysis of a corpus of undeciphered Indus civilization inscriptions indicates syntactic organization. *Computer Speech and Language* **25**, 639-654 (2011).
7. Wells, B. K. *The Archaeology and Epigraphy of Indus Writing* (Archaeopress, Oxford, 2015).
8. Mahadevan, I. *The Indus Script: Texts, Concordance and Tables* (Archaeological Survey of India, Calcutta, 1977).
9. Coulmas, F. *The Blackwell Encyclopedia of Writing Systems* (Blackwell, Malden, MA, 1996).
10. Possehl, G. *The Indus Civilization: A Contemporary Perspective* (AltaMira Press, Lanham, MD, 2002).

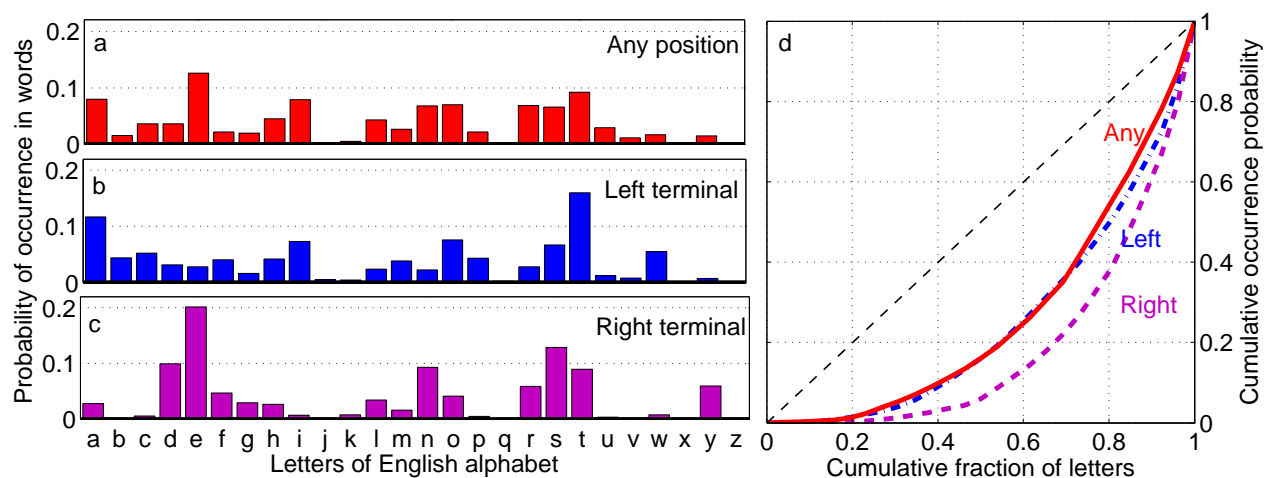


Figure S1. Unequal representation of letters (1-grams) occurring at different positions in words in written English is robust with respect to choice of corpus. The probability of occurrence of the 26 letters of the English alphabet in the *Google Books Ngram* data comprising about 97000 unique words of the English language that occur with a frequency of more than 100,000 in the corpus (see Methods for details), at (a) any position, (b) left terminal position (i.e., in the beginning) and (c) right terminal position (i.e., at the end) of a word. While there are differences in the occurrence probability of the individual letters with the distribution shown in Fig. 1 (see main text), as with the *Mieliestronk* corpus there is higher heterogeneity in (c) indicating that only a few letters occur with high frequency at the right terminal position of a word, compared to a relatively more egalitarian frequency of occurrence of letters in the left terminal position (b). This difference is illustrated in the Lorenz curve (d) comparing the cumulative distribution function for the occurrence probability of the different letters in any (solid curve), left terminal (dash-dotted curve) and right terminal position (dashed curve) of a word. The thin broken diagonal line corresponds to a perfectly uniform distribution, deviation from which indicates the extent of heterogeneity of letter occurrence probability distributions - measured as the ratio of the area between the line of perfect equality and the observed Lorenz curve, i.e., the Gini index.

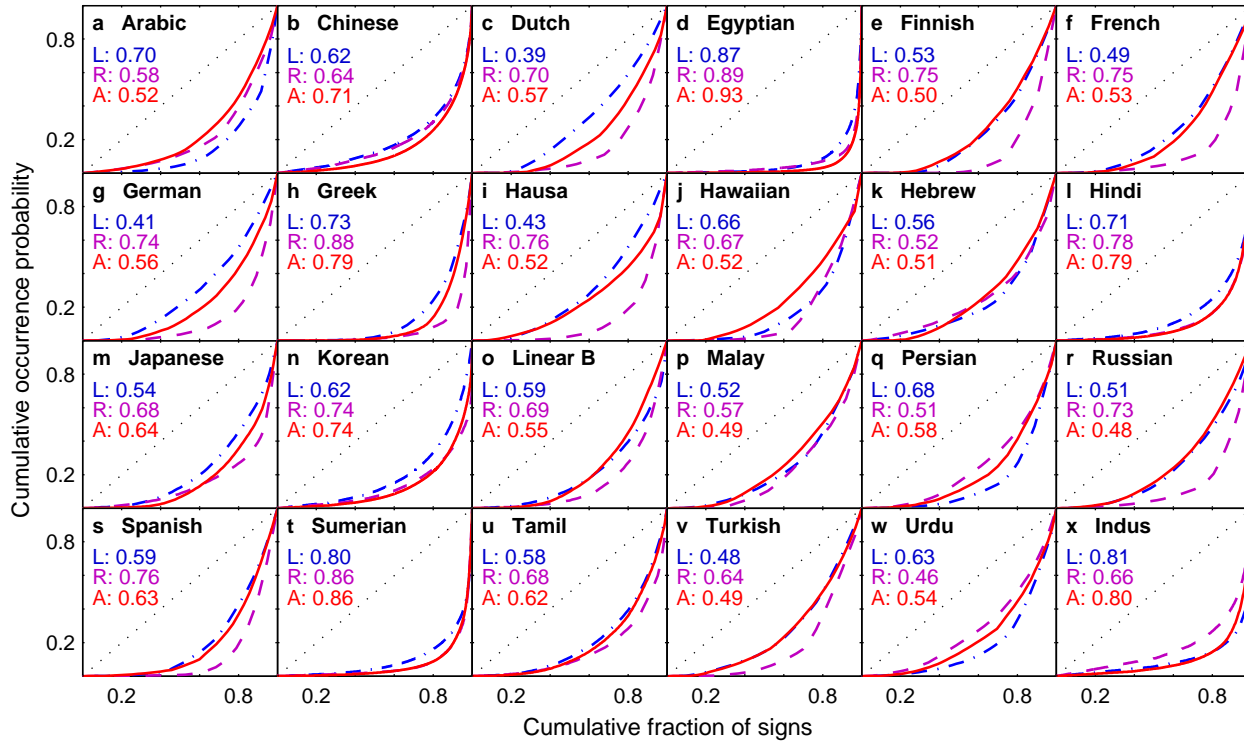


Figure S2. Unequal representation of signs (1-grams) occurring at different positions in words in corpora written using different languages and writing systems. The Lorenz curves in the 24 panels (corresponding to all the scripts analyzed in this paper except English, which is shown in Fig. 1) show the differences in the cumulative distribution function of the occurrence probability of signs at left terminal position (red, dash-dot curve), right terminal position (purple, dashed curve) and at any position (red, solid curve) of a word written in a particular script. The thin broken diagonal line corresponds to a perfectly uniform distribution, deviation from which indicates the extent of heterogeneity of sign occurrence distributions. This is measured in terms of the Gini index (the ratio of the area between the line of perfect equality and the observed Lorenz curve), the corresponding values at the left terminal (L), right terminal (R) and any position (A) for a script being indicated in each panel.

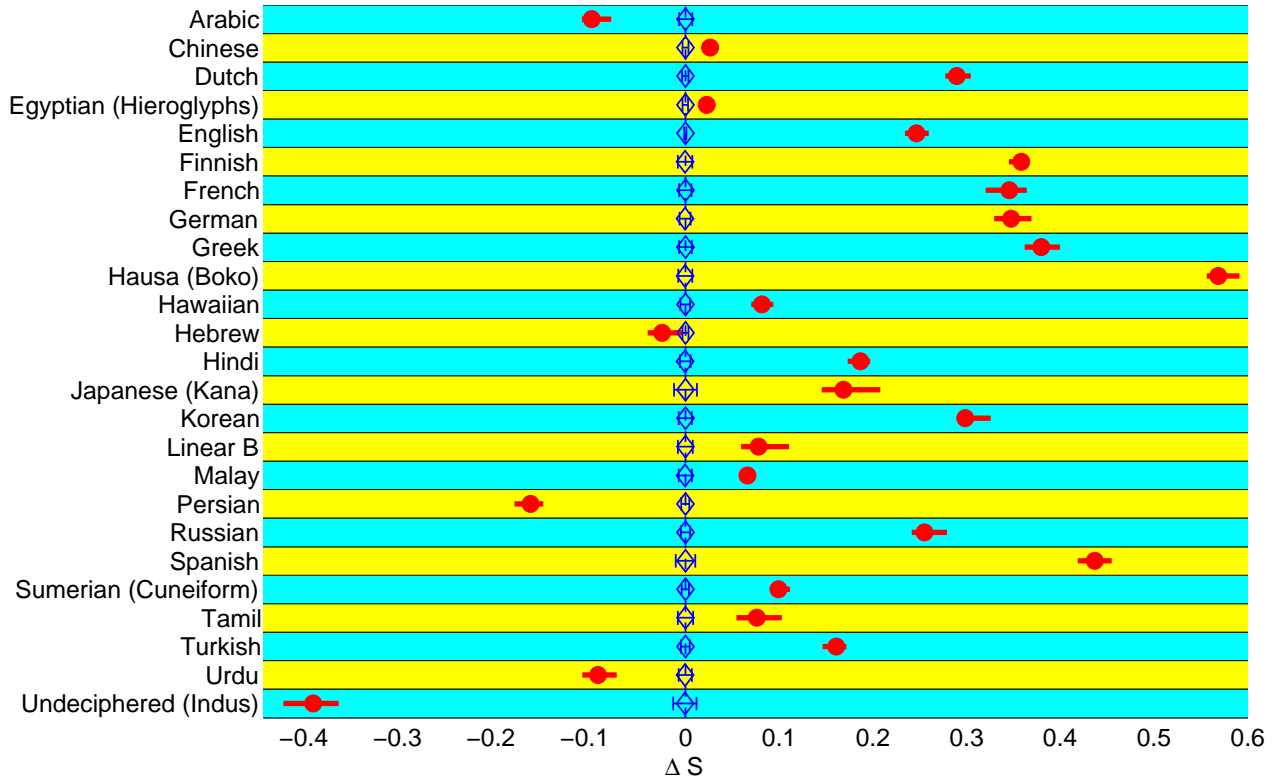


Figure S3. Asymmetry in the sign occurrence probability distributions at the left and right terminal positions of words in different languages are robust with respect to the quantitative measure of inequality used. The normalized difference of the Shannon or information entropies $\Delta S = 2(S_L - S_R)/(S_L + S_R)$ (filled circles), which measures the relative heterogeneity between the occurrences of different signs in the terminal positions of words of a language, are shown for a number of different written languages (arranged in alphabetical order) that span a variety of possible writing systems - from alphabetic (e.g., English) and syllabic (e.g., Japanese kana) to logographic (Chinese) [see text for details]. All languages that are conventionally read from left to right (or rendered in that format in the databases used here) show a positive value for ΔS , while those read right to left exhibit negative values. The horizontal thick bars superposed on the circles represent the bootstrap confidence interval for the estimated values of ΔS . To verify the significance of the empirical values, they are compared with corresponding ΔS (diamonds) calculated using an ensemble of randomized versions for each of the databases (obtained through multiple realizations of random permutations of the signs occurring in each word). Data points are averages over 1000 random realizations, the ranges of fluctuations being indicated by error bars. Along with the set of known languages, ΔS measured for a corpus of undeciphered inscriptions from the Indus Valley Civilization (2600-1900 BCE) is also shown (bottom row).